

## APPENDIX

### A DETAILS OF TENSORS

#### A.1 TENSOR AND MATRIX PRODUCT OPERATORS

As introduced in Gao et al. (2020), a tensor is precisely characterized as follows:

**Tensor.** Let  $D_1, D_2, \dots, D_P \in \mathbb{N}$  denote index upper bounds. A tensor  $\mathcal{T} \in \mathbb{R}^{D_1, D_2, \dots, D_P}$  of order  $P$  is a  $P$ -way array where elements  $\mathcal{T}[d_1, d_2, \dots, d_P]$  are indexed by  $d_p \in \{1, 2, \dots, D_P\}$  for  $1 \leq p \leq P$ .

**Matrix Product Operator.** The concept of bond dimension  $d_k$  is defined as follows:

$$d_k = \min\left(\prod_{p=1}^k i_p \times j_p, \prod_{p=k+1}^n i_p \times j_p\right). \quad (\text{S.1})$$

We can observe that it will be large in the middle and small on both sides by equation S.1. A detailed algorithm for MPO decomposition can be found in Algorithm hdak of Appendix hdakj. The MPO representation decomposes  $M$  into a product of  $n$  local tensors:

$$M_{i_1 i_n, j_1 j_n} = \mathcal{T}^{(1)}[i_1, j_1] \mathcal{T}^{(n)}[i_n, j_n] \quad (\text{S.2})$$

where  $\mathcal{T}^{(p)}[i_p, j_p]$  is a  $D_{p-1} \times D_p$  matrix with  $D_p$  the virtual basis dimension on the bond linking  $\mathcal{T}^{(p)}$  and  $\mathcal{T}^{(p+1)}$  with  $D_0 = D_n = 1$ .

#### A.2 THEOREM

**Theorem 1.** Suppose that the tensor  $\mathbf{W}^{(k)}$  of matrix  $W$  that is satisfy

$$\mathbf{W} = \mathbf{W}^{(k)} + \mathbf{E}^{(k)}, D(\mathbf{W}^{(k)}) = d_k, \text{ where } \|\mathbf{E}^{(k)}\|_F^2 = \epsilon_k^2, k = 1, \dots, d-1. \quad (\text{S.3})$$

Then  $MPO(\mathbf{W})$  with the  $k$ -th bond dimension  $d_k$  upper bound of truncation error satisfy:

$$\|\mathbf{W} - MPO(\mathbf{W})\|_F \leq \sqrt{\sum_{k=1}^{d-1} \epsilon_k^2} \quad (\text{S.4})$$

*Proof.* The proof is by induction. For  $n = 2$  the statement follows from the properties of the SVD. Consider an arbitrary  $n > 2$ . Then the first unfolding  $\mathbf{W}^{(1)}$  is decomposed as:

$$\mathbf{W}^{(1)} = \mathbf{U}_1 \lambda_1 \mathbf{V}_1 + \mathbf{E}^{(1)} = \mathbf{U}_1 \mathbf{B}^{(1)} + \mathbf{E}^{(1)} \quad (\text{S.5})$$

where  $\mathbf{U}_1$  is of size  $r_1 \times i_1 \times j_1$  and  $\|\mathbf{E}^{(1)}\|_F^2 = \epsilon_1^2$ . The matrix  $\mathbf{B}^{(1)}$  is naturally associated with a  $(n-1)$ -dimensional tensor  $\mathcal{B}^{(1)}$  with elements  $\mathcal{B}^{(1)}(\alpha, i_2, j_2, \dots, i_n, j_n)$ , which will be decomposed further. This means that  $\mathbf{B}^{(1)}$  will be approximated by some other matrix  $\hat{\mathbf{B}}^{(1)}$ . From the properties of the SVD it follows that  $\mathbf{U}_1^T \mathbf{E}^{(1)} = 0$ , and thus

$$\begin{aligned} \|\mathbf{W} - \mathcal{B}^{(1)}\|_F^2 &= \|\mathbf{W}_1 - \mathbf{U}_1 \hat{\mathbf{B}}^{(1)}\|_F^2 \\ &= \|\mathbf{W}_1 - \mathbf{U}_1 (\hat{\mathbf{B}}^{(1)} + \mathbf{B}^{(1)} - \mathbf{B}^{(1)})\|_F^2 \\ &= \|\mathbf{W}_1 - \mathbf{U}_1 \mathbf{B}^{(1)}\|_F^2 + \|\mathbf{U}_1 (\hat{\mathbf{B}}^{(1)} - \mathbf{B}^{(1)})\|_F^2 \end{aligned} \quad (\text{S.6})$$

and since  $\mathbf{U}_1$  has orthonormal columns,

$$\|\mathbf{W} - \mathcal{B}^{(1)}\|_F^2 \leq \epsilon_1^2 + \|\mathbf{B}^{(1)} - \hat{\mathbf{B}}^{(1)}\|_F^2. \quad (\text{S.7})$$

and thus it is not difficult to see from the orthonormality of columns of  $\mathbf{U}_1$  that the distance of the  $k$ -th unfolding ( $k = 2, \dots, d_k - 1$ ) of the  $(d-1)$ -dimensional tensor  $\mathcal{B}^{(1)}$  to the  $d_k$ -th rank matrix cannot be larger than  $\epsilon_k$ . Proceeding by induction, we have

$$\|\mathbf{B}^{(1)} - \hat{\mathbf{B}}^{(1)}\|_F^2 \leq \sum_{k=2}^{d-1} \epsilon_k^2, \quad (\text{S.8})$$

combine with Eq. equation S.7, this completes the proof.

## B ALGORITHMS

The MPO pseudocode is shown in Algorithm S.1.

---

**Algorithm S.1** MPO decomposition for a matrix.

---

**Input:** matrix  $\mathbf{M}$ , the number of local tensors  $m$ .  
**Output :** MPO tensor list  $\{\mathcal{T}_{(s)}\}_{s=1}^m$ .

- 1: **for**  $s = 1 \rightarrow m$  **do**
- 2:    $\mathbf{M}[I, J] \rightarrow \mathbf{M}[d_{s-1} \times i_s \times j_s, -1]$
- 3:    $\mathbf{U}\lambda\mathbf{V}^T = \text{SVD}(\mathbf{M})$
- 4:    $\mathbf{U}[d_{s-1} \times i_s \times j_s, d_s] \rightarrow \mathcal{U}[d_{s-1}, i_s, j_s, d_s]$
- 5:    $\mathcal{T}^{(s)} := \mathcal{U}$
- 6:    $\mathbf{M} := \lambda\mathbf{V}^T$
- 7: **end for**
- 8:  $\mathcal{T}^{(s)} := \mathbf{M}$
- 9: Normalization
- 10: **return**  $\{\mathcal{T}_{(k)}\}_{k=1}^n$

---

The over-parameterized matrices selection pseudocode is shown in Algorithm S.2

---

**Algorithm S.2** Fine-tuning a model with our OPF.

---

**Input:** Low rank parameter matrices set of a model  $\{\mathbf{W}\}$ .

- 1: Divide  $\{\mathbf{W}\}$  into several groups by module.
- 2: **if** is Static Strategy **then**
- 3:   LoRA fine-tuning the model until converged.
- 4:   Compute  $I_{\mathbf{W}}$  for  $\{\mathbf{W}\}$  using equation 7.
- 5:   Sort  $\{\mathbf{W}\}$  in each group according to  $I_{\mathbf{W}}$ .
- 6:   Perform MPO on the top- $N$  matrices.
- 7:   Train the other PLM until converged.
- 8: **else**
- 9:   Define  $S = \{\}$
- 10:   **while**  $\text{Len}(S) < N$  **do**
- 11:     Train the model for  $t$  steps.
- 12:     Compute  $I_{\mathbf{W}}$  for  $\{\mathbf{W}\}$  using equation 8.
- 13:     Sort  $\{\mathbf{W}\}$  in each group according to  $I_{\mathbf{W}}$ .
- 14:     Add top- $n$  matrices into  $S$ , and perform MPO.
- 15:   **end while**
- 16:   Continually train the model until converged.
- 17: **end if**

---

## C ADDITIONAL EXPERIMENT DETAILS

In this paper, we propose the MPO decomposition as a method to increase the parameters of the model. Using Eq. 3, an MPO can be specified as:

$$\mathcal{T}_{i_1, i_2, \dots, i_n}^{j_1, j_2, \dots, j_n}(D) \tag{S.9}$$

We pre-compute the significance scores of all parameter matrices before fine-tuning and subsequently over-parameterize the top- $N$  ones using the MPO technique. The significant score can be calculated by Eq. equation 7 and Eq. equation 8.

### C.1 IMPLEMENTATION DETAILS

To ensure a fair comparison, we maintain the experimental setup of GoRA (He et al., 2025) and adopt baseline performances reported by them. By default, we fine-tune the converged models using the AdamW optimizer (Loshchilov & Hutter, 2017) with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e - 8$ . We implement a cosine learning rate schedule with a warmup ratio of 0 and set the rank  $r = 8$  and  $\alpha = 16$ . For natural language understanding tasks, we fine-tune T5-base (Raffel et al., 2020) with a



810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

Datasets	LR	MPO_LR	split number	top- $N$	eval step
<b>Llama 2-7B (Touvron et al., 2023)</b>					
<b>LoRa-Over-SVD</b>					
<b>MT-Bench</b>	1.9e-5	Null	Null	Null	100
<b>GSM8K</b>	1.9e-5	Null	Null	Null	100
<b>HumanEval</b>	1.9e-5	Null	Null	Null	100
<b>LoRa-Over-MPO</b>					
<b>MT-Bench</b>	6e-5	Null	Null	Null	100
<b>GSM8K</b>	6e-5	Null	Null	Null	100
<b>HumanEval</b>	6e-5	Null	Null	Null	100
<b>LoRa-Over-MPO<sub>S</sub></b>					
<b>MT-Bench</b>	6.2e-5	5.6e-5	Null	28	100
<b>GSM8K</b>	6.8e-5	6.1e-5	Null	28	100
<b>HumanEval</b>	6e-5	5.3e-5	Null	28	100
<b>LoRa-Over-MPO<sub>D</sub></b>					
<b>MT-Bench</b>	6.2e-5	5.6e-5	5	28	100
<b>GSM8K</b>	6.8e-5	6.1e-5	7	28	100
<b>HumanEval</b>	6e-5	5.3e-5	5	28	100
<b>Llama 3.1-8B (Grattafiori et al., 2024)</b>					
<b>LoRa-Over-SVD</b>					
<b>MT-Bench</b>	6e-5	Null	Null	Null	50
<b>GSM8K</b>	6e-5	Null	Null	Null	50
<b>HumanEval</b>	6e-5	Null	Null	Null	50
<b>LoRa-Over-MPO</b>					
<b>MT-Bench</b>	6e-5	Null	Null	Null	50
<b>GSM8K</b>	6e-5	Null	Null	Null	50
<b>HumanEval</b>	1e-4	Null	Null	Null	50
<b>LoRa-Over-MPO<sub>S</sub></b>					
<b>MT-Bench</b>	1e-4	9.4e-5	Null	30	50
<b>GSM8K</b>	1e-4	9.5e-5	Null	31	50
<b>HumanEval</b>	1e-4	7e-5	Null	28	50
<b>LoRa-Over-MPO<sub>D</sub></b>					
<b>MT-Bench</b>	1e-4	9.4e-5	5	30	50
<b>GSM8K</b>	1e-4	9.5e-5	5	31	50
<b>HumanEval</b>	1e-4	7e-5	3	28	50

Table S.3: Hyperparameter setup of LoRa-Over for Llama 2-7B and Llama 3.1-8B model. "LR" denote the learning rate. "Null" denote the parameter is useless.

### C.3 EXPERIMENTS ON NATURAL LANGUAGE GENERATION

The hyperparameters of LoRa-Over using Llama 2-7B and Llama 3.1-8B models are presented in the Table S.3. The MPO structure of Llama 2-7B and Llama 3.1-8B model is presented in the Table S.4.

### C.4 COMPARISON OF DIFFERENT LEARNING RATES

Our approach leverages matrix decomposition to over-parameterize low-rank matrices but faces numerical instability, as small perturbations during decomposition can accumulate and compromise model integrity. To address this, we employ MPO decomposition, which ensures near-lossless factorization, enhancing stability and reducing sensitivity to hyperparameter variations. We evaluate robustness by testing learning rates 1.6e-4, 1.8e-4, 2e-4, 2.2e-4, 2.4e-4 on CoLA and MRPC tasks with a T5-Base model (Table S.5). Results show our method is resilient to learning rate changes, with 2e-4 achieving strong performance, minimizing the need for extensive hyperparameter tuning.

